

In [1]:

```
1 import requests
2 from bs4 import BeautifulSoup
3 import pandas as pd
4
```

In [2]:

```
1 url = 'https://www.themoviedb.org/movie?page='
2 page_url = 'https://www.themoviedb.org'
```

In [3]:

```
1 content=requests.get(url,headers={'User-Agent':'Mozilla/5.0'}).text
```

In []:

```
1
```

In [4]:

```
1 ## movie_url of 1st page
2 soup = BeautifulSoup(content,'lxml')
3 movie_url_1st_page = []
4 movie_name_lst = []
5 data = soup.find_all('div',class_ = 'card style_1')
6 for i in data:
7     movie_code = i.a['href']
8     movie_url_1st_page.append(page_url+movie_code)
9     names = i.a['title']
10    movie_name_lst.append(names)
11
```

In [5]:

```
1 #movie_url_1st_page
2
```

In [6]:

```
1 #movie_name_lst
```

```
1 movie_name_lst =[]
2 director_lst = []
3 gen_lst = []
4 run_time = []
5 relese_lst = []
6 raiting_lst = []
7
8 for link in movie_url_1st_page:
9     content=requests.get(link,headers={'User-Agent':'Mozilla/5.0'}).text
10
11    soup = BeautifulSoup(content,'lxml')
```

```
12     data = soup.find_all('div',class_ = 'header_poster_wrapper true')
13
14     raiting = soup.find('div',class_ = 'user_score_chart')['data-percent']
15     raiting_lst.append(raiting)
16
17
18     relese_date = soup.find('span',class_ = 'release')
19
20     director = soup.find('li',class_ = 'profile').a.text
21     director_lst.append(director)
22
23     val = soup.find('span',class_ = 'genres').text.split()
24
25     gen_lst.append(val)
26
27
28     runtime = soup.find('span',class_ = 'runtime').text.split()
29     run_time.append(runtime)
30
31     movie_data_dic = {
32         'Movie Name': movie_name_lst,
33         'Director': director_lst,
34         'Release Date' : relese_lst,
35         'Run Time': run_time,
36         'Genres': gen_lst,
37         'Raiting' : raiting_lst,
38         'Movie_link': movie_url_1st_page
39
40     }
41
42
43
44
45
46
47
```

In []:

1

In []:

1

In []:

1

In []:

1

In []:

```
1  
2
```

In []:

```
1
```

```
1
```

In [7]:

```
1 url_lst = []  
2 for u in range(0,501):  
3     url_lst.append(url+str(u))
```

In []:

```
1
```

In []:

```
1
```

In []:

```
1
```

In []:

```
1
```

In [8]:

```
1 movie_url_all_pages = []  
2 movie_name_lst = []  
3  
4 for link in url_lst:  
5     content=requests.get(link,headers={'User-Agent':'Mozilla/5.0'}).text  
6     soup = BeautifulSoup(content,'lxml')  
7     data = soup.find_all('div',class_ = 'card style_1')  
8  
9     for i in data:  
10         movie_code = i.a['href']  
11         movie_url_all_pages.append(page_url+movie_code)  
12         names = i.a['title']  
13         movie_name_lst.append(names)  
14
```

In []:

1

In [9]:

1 `len(movie_name_lst)`

Out[9]:

10000

In [10]:

```

1 director_lst = []
2 gen_lst = []
3 run_time = []
4 release_lst = []
5 raiting_lst = []
6
7 for link in movie_url_all_pages:
8     content=requests.get(link,headers={'User-Agent':'Mozilla/5.0'}).text
9
10
11     soup = BeautifulSoup(content,'lxml')
12     data = soup.find_all('div',class_ = 'header_poster_wrapper true')
13
14     raiting = soup.find('div',class_ = 'user_score_chart')['data-percent']
15     raiting_lst.append(raiting)
16
17
18     release_date = soup.find('span',class_ = 'release').text.split()[0]
19     release_lst.append(release_date)
20
21     director = soup.find('li',class_ = 'profile')
22     if director is not None:
23         director=(director.p.text)
24     director_lst.append(director)
25
26     val = str(soup.find('span',class_ = 'genres').text)
27     genres = val.replace('\n','')
28     gen_lst.append(genres)
29
30     runtime = soup.find('span',class_='runtime')
31     #if runtime is not None:
32         #runtime=(runtime.text.strip())
33         #run_time.append(runtime)
34     #runtime = (soup.find('span',class_ = 'runtime'))
35     time = runtime.replace('\n','')
36     if runtime is not None:
37         runtime=runtime.text
38
39
40     run_time.append(runtime)
41
42     movie_data_dic = {
43         'Movie Name': movie_name_lst,
44         'Raiting' : raiting_lst,
45         'Release Date' : release_lst,
46         'Run Time': run_time,
47         'Genres': gen_lst,
48         'Director': director_lst,
49         'Movie_link': movie_url_all_pages
50
51     }
52

```

--
TypeError
t)

Traceback (most recent call last)

~\AppData\Local\Temp\ipykernel_11220\193167220.py in <module>

```

33         #run_time.append(runtime)
34         #runtime = (soup.find('span',class_ = 'runtime'))
--> 35         time = runtime.replace('\n','')
36         if runtime is not None:
37             runtime=runtime.text

```

TypeError: 'NoneType' object is not callable

In [12]:

```
1 len(gen_lst)
```

Out[12]:

2536

In [13]:

```
1 df = pd.DataFrame(movie_data_dic)
2 df
```

```

-----
--
ValueError                                Traceback (most recent call last)
~\AppData\Local\Temp\ipykernel_5696\1009747518.py in <module>
----> 1 df = pd.DataFrame(movie_data_dic)
      2 df

~\anaconda3\lib\site-packages\pandas\core\frame.py in __init__(self, data, index, columns, dtype, copy)
    612         elif isinstance(data, dict):
    613             # GH#38939 de facto copy defaults to False only in no
n-dict cases
--> 614         mgr = dict_to_mgr(data, index, columns, dtype=dtype,
copy=copy, typ=manager)
    615         elif isinstance(data, ma.MaskedArray):
    616             import numpy.ma.mrecords as mrecords

~\anaconda3\lib\site-packages\pandas\core\internals\construction.py in dict_to_mgr(data, index, columns, dtype, copy)

```

In []:

```
1
```

In []:

```
1
```

In []:

```
1
```

In []:

1	
---	--

In []:

1	
---	--

In []:

1	
---	--